

(19)



Europäisches Patentamt

European Patent Office

Office européen des brevets



(11)

EP 0 898 236 A2

(12)

EUROPEAN PATENT APPLICATION

(43) Date of publication:

24.02.1999 Bulletin 1999/08

(51) Int. Cl.⁶: G06F 17/30

(21) Application number: 98115643.3

(22) Date of filing: 19.08.1998

(84) Designated Contracting States:

AT BE CH CY DE DK ES FI FR GB GR IE IT LI LU
MC NL PT SE

Designated Extension States:

AL LT LV MK RO SI

(30) Priority: 20.08.1997 JP 223908/97

(71) Applicant:

Toa Gosei Kabushiki Kaisha
Tokyo 105-0003 (JP)

(72) Inventors:

- Yoshida, Tetsuhiko
Nagoya-shi, Aichi 455-0022 (JP)
- Osawa, Kenji
Nara-shi, Nara 631-0011 (JP)
- Obata, Nobuaki
Nagoya-shi, Aichi 464-0096 (JP)

(74) Representative:

R.A. KUHNEN & P.A. WACKER
Patentanwalts-gesellschaft mbH
Alois-Steinecker-Strasse 22
85354 Freising (DE)

(54) Method and apparatus for manifesting characteristic existing in symbolic sequence

(57) A method which manifests characteristic which is latent and can not be recognized, although it exists in a complicated symbolic sequence, for example, a nucleotide sequence of DNA, and thereby enables recognition of the characteristic unrecognized yet, is provided.

When a symbolic sequence I_j ($j = 1 \sim m$) is given, there is an effected conversion to a parallel sequence $A(k)$ of partial symbolic sequences in which the suffix j is aligned in the following positional relation:

$$\begin{array}{ccccccc}
 j = & 1, & 2, & \dots & k-1, & k \\
 j = & k+1, & k+2, & \dots & k+k-1, & k+k \\
 : & & & & & \\
 : & & & & & \\
 J = & (n-1)k+1, & (n-1)k+2, & \dots & (n-1)k+k-1, & (n-1)k+k \\
 j = & nk+1, & nk+2, & \dots & nk+k-1, & nk+k
 \end{array}$$

and $A(k)$ is formed with changing k to $p, p+r, p+2r, p+3r \dots$, and the whole parallel sequences $\Sigma A(k)$ is obtained.

When regularity of period length k exists in the symbolic sequence I_j , the regularity remarkably appears in the partial symbolic sequences obtained by extracting one symbol at every $k-1$ symbols from the symbolic sequence.

EP 0 898 236 A2

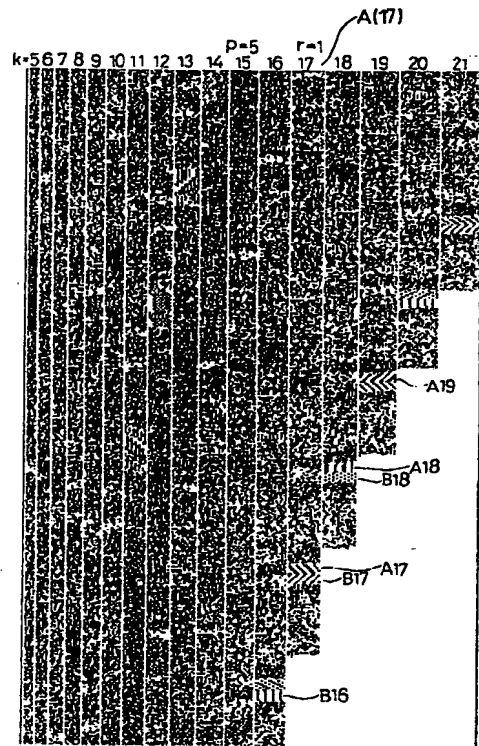


FIG.1

Description

BACKGROUND OF THE INVENTION

5 Field of the Invention

[0001] The present invention relates to a method and apparatus for manifesting characteristic or regularity which is latent and can not be recognized, although such characteristic or regularity actually exists in a complicated symbolic sequence, for example, a nucleotide sequence of DNA, an amino acid sequence of protein, or a digital sequence of decimal expansion of an irrational number and the like. In these sequences, regularity can not be recognized at a glance even when the regularity is included. The present invention enables recognition of characteristic or regularity included in the symbolic sequence but unrecognized yet.

Description of the Related Art

15

[0002] Some complicated symbolic sequences contain characteristic which has not been recognized by human beings, although the characteristic actually exists. For example, genetic information is specified by a symbolic sequence of long string. The symbols consist of four symbols each indicating one of four kinds of nucleotides. A large amount of symbols are one-dimensionally aligned. In the study of genetic information, it is extremely important to recognize a certain regularity hidden in the symbolic sequence indicating genetic information. Besides, if a certain regularity is found in an irrational number, the number π , and the base of natural logarithm (e), the study of random numbers is intensified and various developments are expected in mathematics.

[0003] For such a purpose, various trials have been made for analyzing a symbolic sequence based on a variety of mathematical methods such as the Fourier analysis. However, these trials have not necessarily accomplished successful results. One problem with the conventional analysis methods is that even if a certain regularity is included in a part of a very long symbolic sequence, the regularity existing partially is buried in the whole sequence and can not be recognized when the whole symbolic sequence is analyzed. Since there is no effective technology to know in advance in which part of the sequence the regularity exists, there are many characteristics or regularity which can not be recognized by the conventional analysis methods.

30

SUMMARY OF THE INVENTION

[0004] An object of the present invention is to contrive a method and apparatus which manifests characteristic or regularity even if such characteristic or regularity exists only in a part of the whole symbolic sequence, and thereby enables recognition of characteristic or regularity which has not been recognized until now.

[0005] Another object of the present invention is to manifest characteristic or regularity existing throughout the entire sequence.

[0006] In one embodiment of the present invention, when a symbolic sequence I_j ($j = 1 \sim m$) is given, there is an effected conversion to a parallel sequence $A(k)$ of partial symbolic sequences in which the suffix j is aligned in the following positional relation:

$$\begin{array}{rcll}
 j = & 1, & 2, \dots\dots\dots & k-1, & k \\
 j = & k+1, & k+2, \dots\dots\dots & k+k-1, & k+k \\
 : & & & & \\
 : & & & & \\
 : & & & & \\
 J = & (n-1)k+1, & (n-1)k+2, \dots\dots\dots & (n-1)k+k-1, & (n-1)k+k \\
 j = & nk+1, & nk+2, \dots\dots\dots & nk+k-1, & nk+k.
 \end{array}$$

Instead, the positional relation may be the following :

$$\begin{array}{l}
 j = 1, \quad 2, \dots, \quad k-1, \quad k \\
 j = k+k, \quad k+k-1, \dots, \quad k+2, \quad k+1 \\
 : \\
 : \\
 j = (n-1)k+k, \quad (n-1)k+k-1, \dots, \quad (n-1)k+2, \quad (n-1)k+1 \\
 j = nk+1, \quad nk+2, \dots, \quad nk+k-1, \quad nk+k.
 \end{array}$$

Herein, k represents an integer of 2 or more, n represents an integer such that $nk < m \leq nk+k$, and when the suffix j is $m+1$ or more, the processing is ignored.

[0007] Then, the converted parallel sequence $A(k)$ is output using one or more expression means selected from hue, lightness and saturation of color and from interval, tone and volume of sound.

[0008] Equidistant Letter Sequences in the Book of Genesis (Doron Witztum, Eliyahu Rips and Yoav Rosenberg, Statistical Science 1994, Vol. 9, No. 3, page 429-438) introduces a technology in which a code hidden in a one-dimensional letter sequence is decoded by converting the one-dimensional letter sequence to a parallel sequence $A(k)$ of partial symbolic sequences. In this technology, there is required an operation to extract letters having senses from the parallel sequence $A(k)$ of partial symbolic sequences, and it can not be used for other sequences than the letter sequence. Further, when a certain regularity is to be recognized in a symbolic sequence which is irregular at a glance and which is often found in a natural field, inconsistency, namely, regularity can not be recognized unless the regularity to be recognized has been recognized in advance, can not be solved.

[0009] In the present invention described above, since a parallel sequence $A(k)$ of partial symbolic sequences is output using one or more expression means selected from hue, lightness and saturation of color and from interval, tone and volume of sound, even if regularity is not known in advance, that regularity is manifested by a pattern of hue, lightness and saturation of color and interval, tone and volume of sound and can be easily recognized.

[0010] In another embodiment of the present invention, when a one-dimensional symbolic sequence I_j ($j = 1 \sim m$) is given, there is an effected conversion to a parallel sequence $A(k)$ of partial symbolic sequences in which the suffix j is aligned in the following positional relation:

$$\begin{array}{l}
 j = 1, \quad 2, \dots, \quad k-1, \quad k \\
 j = k+1, \quad k+2, \dots, \quad k+k-1, \quad k+k \\
 : \\
 : \\
 j = (n-1)k+1, \quad (n-1)k+2, \dots, \quad (n-1)k+k-1, \quad (n-1)k+k \\
 j = nk+1, \quad nk+2, \dots, \quad nk+k-1, \quad nk+k.
 \end{array}$$

Instead, the positional relation may be the following :

$$\begin{array}{llll}
 j = & 1, & 2, & \dots\dots\dots k-1, & k \\
 j = & k+k, & k+k-1, & \dots\dots\dots k+2, & k+1 \\
 & : & & & \\
 & : & & & \\
 & : & & & \\
 j = & (n-1)k+k, & (n-1)k+k-1, & \dots\dots\dots (n-1)k+2, & (n-1)k+1 \\
 j = & nk+1, & nk+2, & \dots\dots\dots nk+k-1, & nk+k.
 \end{array}$$

Further, when p represents a natural number from 2 to less than m , r represents any natural number, the above-described conversion is repeated with changing k to $p, p+r, p+2r, p+3r, \dots\dots\dots$ to obtain parallel sequences of partial symbolic sequences: $A(p), A(p+r), A(p+2r), A(p+3r), \dots\dots\dots$. Then, the resulted parallel sequences: $A(p), A(p+r), A(p+2r), A(p+3r), \dots\dots\dots$ are further parallel-positioned to make a whole parallel sequences $\Sigma A(k)$. Then, the obtained whole parallel sequences $\Sigma A(k)$ is output. Here, n represents an integer such that $nk < m \leq nk+k$, and when the suffix j is $m+1$ or more, the processing is ignored.

[0011] In this case, a parallel sequence made by parallel-positioning of p partial symbolic sequences, a parallel sequence made by parallel-positioning of $p+r$ partial symbolic sequences, and parallel sequences made by parallel-positioning of likewise increased number of partial symbolic sequences, are all parallel-positioned. In this processing, if regularity of period length α is hidden in the symbolic sequence, such regularity is remarkably manifested in a parallel sequence $A(\alpha)$ made by parallel-positioning of partial symbolic sequences of a number of α .

[0012] If α is included in between $p, p+r, p+2r, p+3r, \dots\dots\dots$ rows, regularity of period length α is manifested in a parallel sequence of partial symbolic sequences of an analogous number to α . Therefore, increment r regarding the number of the partial symbolic sequences is not necessarily required to be one, and it may advantageously be any natural number. In this case, when the increment r is smaller, characteristic is more securely manifested.

[0013] In this embodiment, regularity of an unknown period length is manifested in a parallel sequence of partial symbolic sequences of some number, and recognition of characteristic becomes easy.

[0014] In the above-described method, it is preferable that each symbol is expressed by combination of hue, lightness and saturation of color. By this embodiment, as a result of the manifestation of characteristic hidden in the symbolic sequence through visual sense, there will be made more sufficient understanding regarding characteristic hidden in the symbolic sequence, and various applications and developments utilizing the characteristic are made possible. Further, the resulting visual pattern is a pattern including regularity and irregularity mixed which has not conventionally existed, and a visual pattern of which design itself has utility application can be designed.

[0015] Each symbol may be expressed by combination of interval, tone and volume of sound. By expressing the parallel sequence by combination of interval, tone and volume of sound, a unique audio pattern is created and the characteristic of the symbolic sequence can be recognized through auditory sense.

[0016] When one symbol is taken out from an original symbolic sequence at an interval of $k-1$ (namely, at every k) to make a symbolic sequence and the method is applied to this extracted symbolic sequence, if regularity of period length k is hidden in the original symbolic sequence, the regularity is manifested and appears remarkably.

[0017] When one symbol is taken out from an original symbolic sequence at an interval of $kq-1$ (namely, at every kq) to make a symbolic sequence and the method is applied to this extracted symbolic sequence, if regularity of period length kq is hidden in the original symbolic sequence, the regularity is manifested and appears remarkably.

[0018] Further, when any of the above-described methods are conducted with changing k to $p, p+r, p+2r, \dots\dots\dots$, there is formed a whole parallel sequences $\Sigma A(k)$ of a parallel sequence $A(p)$ made by parallel-positioning of p partial symbolic sequences, a parallel sequence $A(p+r)$ made by parallel-positioning of $p+r$ partial symbolic sequences, and parallel sequences made by parallel-positioning of likewise increased number of partial symbolic sequences, and regularity of period length α appears remarkably in a parallel sequence $A(\alpha)$ formed by parallel-positioning of $k (= \alpha)$ partial symbolic sequences. Therefore, characteristic or regularity of unknown period length is manifested, and recognition of characteristic or regularity becomes easy.

[0019] According to this method, even if the period length α of regularity or characteristic is included between $p, p+r, p+2r, \dots\dots\dots$ rows, characteristic or regularity is manifested in a parallel sequence formed by parallel-positioning partial

sequences of a number approximated to α , and increment r is not necessarily required to be one. It becomes possible to manifest characteristic by small amounts of data processing, by selecting increment r adjusted to events.

[0020] Further, by outputting analyzed results by color and/or sound, expressions suitable to event and observer become possible, and characteristic is more easily recognizable. The resulted color and/or sound pattern will be an interesting pattern in which regularity and irregularity are mixed, and the method can be utilized also as a designing method.

[0021] Especially when an initiation position of regularity is situated at an analysis initiation position, the pattern of having a parabolic shape clearly appears, and regularity of a long period length is manifested clearly, in the whole parallel sequences $\Sigma A(k)$.

[0022] The present invention will be recognized more successfully by reading the descriptions of the following examples with the referring drawings.

BRIEF EXPLANATION OF THE DRAWINGS

[0023]

Fig. 1 represents the whole parallel sequences $\Sigma A(k)$ formed from a nucleotide sequence of human genomic DNA.

Fig. 2 represents positional relation of suffix j in Fig. 1.

Fig. 3 represents the whole parallel sequences $\Sigma A(k)$ formed from a numerical sequence showing π .

Fig. 4 represents the whole parallel sequences $\Sigma A(k)$ formed from a circulating numerical sequence of period length 18 (symbolic sequence).

Fig. 5 represents the whole parallel sequences $\Sigma A(k)$ formed from a circulating numerical sequence of period length 12 (symbolic sequence).

Fig. 6 represents a part of the whole parallel sequences $\Sigma A(k)$ formed from an amino acid sequence of muscle protein myosin.

Fig. 7 represents another part of the whole parallel sequences $\Sigma A(k)$ formed from the amino acid sequence of muscle protein myosin.

Fig. 8 represents still another part of the whole parallel sequences $\Sigma A(k)$ formed from the amino acid sequence of muscle protein myosin.

Fig. 9 represents existing positions of vowel 'O' appearing in 'Genji Monogatari'.

Fig. 10 represents placing positions in converting 100 symbolic sequence into the whole parallel sequences $\Sigma A(k)$.

Fig. 11 explains pre-treatments for symbolic sequences.

Fig. 12 represents another example of placing positions in converting a symbolic sequence into the whole parallel sequences $\Sigma A(k)$.

Fig. 13 represents the whole parallel sequences $\Sigma A(k)$ formed from cDNA sequence of a G protein β subunit.

Fig. 14 represents extraction of symbolic sequence I to be processed with changing the initial point, from a symbolic sequence M .

Fig. 15 represents an example in which parabolic pattern appears in the whole parallel sequences $\Sigma A(k)$ formed from a genomic DNA sequence of baker's yeast.

Fig. 16 represents another positioning (reciprocal) pattern for obtaining a parallel sequence $A(k)$.

Fig. 17 represents the whole parallel sequences $\Sigma A(k)$ formed from a circulating sequence of period length 100.

Fig. 18 represents constitution of an apparatus for effecting the method of the present invention.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

[0024] Experimental examples embodying the present invention will be introduced below.

[First experimental example]

[0025] Fig. 1 represents an experimental example for processing a symbolic sequence I_j showing nucleotide sequence of human genomic DNA. A symbolic sequence I_j indicating nucleotide sequence of human genomic DNA is formed by one-dimensional sequence of an enormous amount of symbols, each symbol indicating one of four kinds of nucleotides ATGC, and a certain regularity hidden therein is recognized as useful information. Therefore, it is a large object in genetic study to find regularity, or to specify a part of the sequence including the regularity.

[0026] Fig. 1 represents a processed result output through color, and four kinds of symbols ATGC are, respectively, expressed by four colors red, blue, green and yellow. The original image of Fig. 1 is expressed in four colors. Fig. 1 represents a result when the present invention is conducted with $p=5$ and $r=1$.

[0027] Referring to an example of a parallel sequence $A(17)$ in which $k=17$, as shown in Fig. 2, longitudinal partial

symbolic sequences C1, C2, C3 ... C17 which are extracted from a symbolic sequence I_j at every k and aligned longitudinally, are laterally aligned to form a parallel sequence A(17). In C1, C2, C3 ..., values of suffixes j of symbols to be extracted are shifted by one. This rule is common to all k values and to all partial symbolic sequences C.

[0028] In this example, a symbol group extracted at every k is placed longitudinally to form a longitudinal partial symbolic sequence, and the longitudinal partial symbolic sequences are laterally placed. However, longitudinal to lateral relations may be reversed, and a symbol group extracted at every k may be placed laterally to form a lateral partial symbolic sequence, and the lateral partial symbolic sequences may be longitudinally placed.

[0029] In Fig. 1, B16 remarkably shows that a repeating pattern of period length 16 exists in a part of the nucleotide sequence. From the pattern B16, it is possible to learn that there is a possibility that useful information is included in this part, and this part is a area valuable to be analyzed in detail. B17 and B16 represent the same regularity. The regularity of period length 16 appears as the vertical stripes in B16, and appears as the inclined stripes in B17. The inclined stripes in B17 has a pattern in which the left side is lowered. B18 also represents the same regularity, and the inclination of the stripes in B18 is closer to horizontal than in B17. The same regularity is also shown in a parallel sequence A(19) in which $k=19$, however, in this case, the inclination is almost horizontal, and extraction of characteristic becomes increasingly difficult.

[0030] Regularity of period length α appears vertically and is expressed most remarkably in A(α) in which $k(=\alpha)$ partial symbolic sequences are parallel-positioned. However, the regularity also appears in a parallel sequence of partial symbolic sequences in which $k=\alpha+1$ and $k=\alpha+2$. Therefore, it is confirmed that the increment r is not necessarily required to be 1.

[0031] A18 shows regularity of period length 18, and the same regularity is shown as pattern A17 in which the right side is lowered in the parallel sequence A(17) in which $k=17$, and shown as pattern A19 in which the left side is lowered in the parallel sequence A(19) in which $k=19$.

[0032] In addition, many remarkable patterns appear in Fig. 1, and characteristics hidden in a nucleotide sequence of human genomic DNA are grasped from these patterns.

[0033] Initial number p in the whole parallel sequences of partial symbolic sequences may be any natural number, and in Fig. 1, $p=5$. The increment r is not limited to 1, and it may be 2 or more. When r is smaller, characteristics are never failed to be found, and when r is larger, the amount of data processing is smaller. The increment r is not required to be constant, and it is preferable to select the increment r according to events.

[0034] Fig. 18 represents an apparatus for effecting the above-described processing method, and in this apparatus, a symbolic sequence I_j to be analyzed is memorized in a memory apparatus 181, and an apparatus 182 converts the symbolic sequence I_j into a parallel sequence A(k), and apparatus 183 forms the whole parallel sequences $\Sigma A(k)$ in which a plurality of parallel sequences A(k) obtained by changing the value of k are parallel-positioned, and an apparatus 184 outputs the whole parallel sequences $\Sigma A(k)$. The apparatus 182 and the apparatus 183 may be constituted of a computer and the apparatus 184 may be constituted of a color printer. When the whole parallel sequences $\Sigma A(k)$ is output using sound, a sound synthesizer may be used as the apparatus 184.

[0035] It is preferable that Fig. 1 is expressed with a time lapse according to speed for processing the symbolic sequence. For example in Fig. 1, color corresponding to I1 is first expressed on the left upper summits of A(5) to A(21), and further expressions of I2, I3, I4 ... are effected in succession. By using this change in time, characteristics are more easily recognized, and also in the case of output by sound, output with time lapse is effective, and when output with time lapse is conducted, characteristics are recognized through changes in sound.

[0036] Fig. 3 exemplifies a result obtained by processing π symbolic sequence (numerical sequence), and 10 kinds of symbols (number) 0 to 9 are expressed by 10 equally divided colors of spectrum from a violet to red. It is found from the expression result of Fig. 3 that specific symbols (numbers) tend to appear frequently in a specific range.

[0037] When noise input is processed as a row of a symbolic sequence and this symbolic sequence is processed to obtain similar pattern as in Fig. 3, it becomes possible to extract characteristic existing in the noise and to extract only meaningful sound included in the noise. Further, it is known that the pattern shown in Fig 3 can be used, for example, as a ground pattern for securities, and this complicated ground pattern can be specified by a one-dimensional symbolic sequence.

[0038] Fig. 4 represents a result obtained by processing a circulating numerical sequence of period length 18, and various patterns can be drawn according to the number k of a partial symbolic sequences to be fractionated. Various textile patterns can be designed by this pattern creating technology. Fig. 5 represents a processed result of a circulating numerical sequence of period length 12, and it is confirmed that different patterns from those of Fig. 4 can be made. According to this method, the complicated pattern shown in Fig. 3 and the regular patterns shown in Figs. 4 and 5 can be designed by the same method. Further, various patterns having utterly different impressions can be produced by changing corresponding relations of a symbol to a color.

[0039] Figs. 6 through 8 represent a result obtained by processing a symbolic sequence which shows an amino acid sequence of a protein myosin of an adductor muscle of a scallop. In Figs. 6 to 8, basic residue is shown in blue, polar residue is shown in green, acidic residue is shown in red, and hydrophobic residue is shown in yellow. In Fig. 6, a

remarkable yellow longitudinal stripe appears in a parallel sequence in which $k=7$, and the existence of regularity of α -helix of period length 7 is found. This regularity of a hydrophobic residue of period length 7 corresponds to α -helix, and by this method, the existence of α -helix can be recognized and the existing position thereof can be specified. This α -helix is manifested as yellow longitudinal stripes in parallel sequences in which $k=7, 14, 28$ and 35 , and manifested as yellow diagonal lines in parallel sequences in which, for example, $k=22, 27$ and 29 .

[0040] Fig. 9 represents an example expressing dots in positions where vowel 'O' appears, prepared by applying the present invention to a symbolic sequence showing a row of vowels in Genji Monogatari. The left side represents an analysis result of the Kiritsubo chapter, and the right side represents an analysis result of the Hahakigi chapter. There is manifested characteristic that appear with frequency in vowel 'O' as high in the specific range of the document, and as low in other specific ranges. By this method, extraction of characteristic in a letter information becomes easy.

[0041] Fig. 10 schematically represents processing contents for a symbolic sequence I_j ($j=1$ to 100) to be processed.

[0042] When, the period length of regularity to be extracted is known in advance, it will be easily recognized whether the regularity of the known period length k really exists, and in the case of existence, where it exists, by forming a parallel sequence $A(k)$ in which partial symbolic sequences obtained by division into k fractions are parallel-positioned.

[0043] Even when the period length is not known, regularity of the unknown period length is manifested at some location in the whole parallel sequences.

[0044] Fig. 11 represents an example of pre-processing for a symbolic sequence to be processed, and when part of a symbolic sequence J shown in (A) is processed, the part to be processed as shown in (B) will be the whole symbolic sequence I of the present invention. Further, when one symbol is specified by combination of a plurality of symbols, this

method is applied for the symbolic sequence specified by combination of a plurality of symbols, for example, as shown in (C). Alternatively, it may also be permissible that one symbol is obtained from symbols of order 123 in a symbolic sequence K , then, one symbol is obtained from symbols of order 234 in a symbolic sequence K , this procedure is repeated to effect conversion into one symbolic sequence I , and this converted symbol I is processed by the method, as in the case for calculating moving average. Further, as shown in (E), for a symbolic sequence existing in a symbolic sequence at specific period, it may also be permissible that a symbolic sequence of this period is first extracted, and the present invention is applied to the extracted symbolic sequence.

[0045] Instead of this method, processing as exemplified in Fig. 12 may be effected. In this method, one symbol is extracted at every kq , for a partial symbolic sequence of longitudinal direction. In the case shown in this drawing, the result is obtained by effecting the method with changing k to $2, 3, 4 \dots$ and q is fixed at 5 . The result corresponds to the result when a symbol of an order of $5 \cdot 10 \cdot 15 \dots$ is first extracted, and then the extracted sequence is separated into k partial symbolic sequences, and the resulted partial sequences are parallel-positioned to obtain a parallel sequence. By this method, it is possible to manifest regularity further hidden in a symbolic sequence hidden in a symbolic sequence L (shown in (E) of Fig. 11).

[0046] When a parallel sequence of partial symbolic sequences are obtained as described above, various method can be adopted for expressing the result, and a method in which a symbol is expressed by color, a method in which a symbol is expressed by variation in color density and a method in which a symbol is expressed by a character (two dimensional pattern) may be adopted, and further, the resulted line and row of symbols may also be expressed by sound. In this case, chord is made by an arrangement of symbols in line direction, and an arrangement in a row direction is expressed by changing this chord by time. By this procedure, it becomes possible to grasp characteristic existing in a symbolic sequence through sound.

[0047] The present invention is useful for analyzing various symbolic sequences, and useful in analyzing a nucleotide sequence of DNA, a nucleotide sequence of RNA, an amino acid sequence of protein, a numerical sequence, a letter sequence, a sound sequence and the like. By this analysis, it becomes possible to specify an existing position of useful information and to extract useful information. Further, when this method is applied to two symbolic sequences which can not be distinguished at a glance, characteristics are manifested, and the identity can be easily judged. In this sense, characteristics and regularity manifested in this method are not restricted to a repeating pattern having a certain period, and characteristics found in distribution of appearing sequence are also manifested. Further, the increment r in the number of partial symbolic sequences is not necessarily required to be 1 , and further, it may not be a constant number. By effecting this method according to $k_1, k_2, k_3 \dots$ distributing irregularly, characteristics existing in two or more symbolic sequences are manifested, and the identity is easily judged.

[0048] Fig. 13 represents the analysis result of a cDNA sequence of a G protein β subunit, and represents the result when the whole parallel sequences $\Sigma A(k)$ is obtained when p is 5 . In the original image of Fig. 13, GCTA are expressed by 4 colors and three apparent different color zones are recognized.

[0049] The boundary 101 of the color zones corresponds approximately to the position of $j=281$, and the boundary 102 of the color zones corresponds approximately to the position of $j=1303$. In this case, it is known that a coding range exists in the range from $j=281$ to $j=1303$, and it is recognized that the coding range is easily specified through visual sense in this method.

[0050] Fig. 14 represents a procedure to obtain symbolic sequence I to be processed when changing the initial point,

from a one-dimensional symbolic sequence M. For example, symbolic sequence I6 to be processed is a symbolic sequence obtained by extraction of M(6) and the following.

[0051] When the present invention is performed on symbolic sequences I1, I2, I3, I4 ... thus extracted to obtain the whole parallel sequences $\Sigma A(k)$, a clear pattern may be obtained in the whole parallel sequences $\Sigma A(k)$ corresponding to specific I.

[0052] Fig. 15 represents one example thereof, and in the original image which is expressed in multicolor, a plurality of parabolic lines 151, 152, 153 ... appear.

[0053] As a result of intensive study of this phenomenon, it has been recognized that the above-described line group appears when the initiation point of regularity coincides with the initiation point of the symbolic sequence to be processed. By this, it has been known that the initiation point of regularity can be specified by utilizing appearance of a line group.

[0054] Further, it is also known that the appearance gap of a group of lines 151, 152, 153 ... and other group of lines 161, 162, 163 ... corresponds to regularity of an extremely long period, and it has also been recognized that the regularity of an extremely long period can be recognized by utilizing a line group.

[0055] It has been recognized that the above-described pattern appears also by reversing alternately a sequential direction of partial symbolic sequences in a lateral direction (reciprocal positioning pattern). Fig. 16 represents a positional relation for obtaining a parallel sequence A(k) by reversing alternately a sequential direction of partial symbolic sequences in lateral direction. Fig. 17 represents an example in which a circulating sequence having a period of 100 is converted to the whole parallel sequences $\Sigma A(k)$ having the positional relation as shown in Fig. 16, and it is recognized that a clear line group appears.

[0056] The above-described explanations are only some specific examples and the present invention can be used in various ways within the attached claims.

Claims

1. A method for manifesting characteristic existing in a symbolic sequence I_j ($j = 1 \sim m$), comprising the steps of:

effecting conversion to a parallel sequence A(k) of partial symbolic sequences in which the suffix j is aligned in the following positional relation:

$j = 1, 2, \dots, k-1, k$

$j = k+1, k+2, \dots, k+k-1, k+k$

:

:

:

$J = (n-1)k+1, (n-1)k+2, \dots, (n-1)k+k-1, (n-1)k+k$

$j = nk+1, nk+2, \dots, nk+k-1, nk+k$

or

$j = 1, 2, \dots, k-1, k$

$j = k+k, k+k-1, \dots, k+2, k+1$

:

:

:

$j = (n-1)k+k, (n-1)k+k-1, \dots, (n-1)k+2, (n-1)k+1$

$j = nk+1, nk+2, \dots, nk+k-1, nk+k;$

and

outputting the converted parallel sequence $A(k)$ using one or more expression means selected from hue, lightness and saturation of color and from interval, tone and volume of sound;

wherein, k represents an integer of 2 or more, and n represents an integer such that $nk < m \leq nk+k$, and when the suffix j is $m+1$ or more, the processing is ignored.

2. A method for manifesting characteristic existing in a symbolic sequence I_j ($j = 1 \sim m$), comprising the steps of:

effecting conversion to a parallel sequence $A(k)$ of partial symbolic sequences in which the suffix j is aligned in the following positional relation:

$j = 1, 2, \dots, k-1, k$

$j = k+1, k+2, \dots, k+k-1, k+k$

:

:

:

$J = (n-1)k+1, (n-1)k+2, \dots, (n-1)k+k-1, (n-1)k+k$

$j = nk+1, nk+2, \dots, nk+k-1, nk+k$

or

$j = 1, 2, \dots, k-1, k$

$j = k+k, k+k-1, \dots, k+2, k+1$

:

:

:

$j = (n-1)k+k, (n-1)k+k-1, \dots, (n-1)k+2, (n-1)k+1$

$j = nk+1, nk+2, \dots, nk+k-1, nk+k ;$

making a whole parallel sequences $\Sigma A(k)$ by further parallel-positioning of parallel sequences $A(p), A(p+r), A(p+2r), A(p+3r), \dots$ converted by changing k to $p, p+r, p+2r, p+3r, \dots$, wherein p represents a natural number from 2 to less than m , and r represents any natural number; and

outputting the obtained whole parallel sequences $\Sigma A(k)$, wherein n represents an integer such that $nk < m \leq nk+k$, and when the suffix j is $m+1$ or more, the processing is ignored.

3. The method according to Claim 2, wherein the whole parallel sequences $\Sigma A(k)$ is output by using one or more expression means selected from hue, lightness and saturation of color and from interval, tone and volume of sound.

4. The method according to Claim 2, further comprising the steps of:

making the symbolic sequence I_j by extracting symbols sequentially from a symbolic sequence M_s ($s=1$ to u and $u > m$);

making the whole parallel sequences $\Sigma A(k)$ from the extracted symbolic sequence I_j ; and

repeating said two steps with shifting an initiation point at which the symbolic sequence I_j extracted from the symbolic sequence M_s .

5. The method according to Claim 2, further comprising the step of:

making the symbolic sequence I_j by taking out m symbols sequentially at every q from a symbolic sequence L_s ($s=1$ to t and $t > mq$).

6. The method according to Claim 1, further comprising the step of:

making the symbolic sequence I_j by taking out m symbols sequentially at every q from a symbolic sequence L_s ($s=1$ to t and $t>mq$).

7. An apparatus for manifesting characteristic existing in a symbolic sequence I_j ($j = 1 \sim m$), comprising:

a means for memorizing the symbolic sequence I_j ;

a means for effecting conversion to a parallel sequence $A(k)$ of partial symbolic sequences in which the suffix j is aligned in the following positional relation:

$$j = \quad 1, \quad 2, \dots \quad k-1, \quad k$$

$$j = \quad k+1, \quad k+2, \dots \quad k+k-1, \quad k+k$$

:

:

:

$$J = (n-1)k+1, (n-1)k+2, \dots (n-1)k+k-1, (n-1)k+k$$

$$j = \quad nk+1, \quad nk+2, \dots \quad nk+k-1, \quad nk+k$$

or

$$j = \quad 1, \quad 2, \dots \quad k-1, \quad k$$

$$j = \quad k+k, \quad k+k-1, \dots \quad k+2, \quad k+1$$

:

:

$$j = (n-1)k+k, (n-1)k+k-1, \dots (n-1)k+2, (n-1)k+1$$

$$j = \quad nk+1, \quad nk+2, \dots \quad nk+k-1, \quad nk+k;$$

a means for making a whole parallel sequences $\Sigma A(k)$ by further parallel-positioning of parallel sequences $A(p)$, $A(p+r)$, $A(p+2r)$, $A(p+3r)$... converted by changing k to p , $p+r$, $p+2r$, $p+3r$, ... when p represents a natural number from 2 to less than m , and r represents an any natural number; and

a means for outputting the obtained whole parallel sequences $\Sigma A(k)$, wherein n represents an integer such that $nk < m \leq nk+k$, and when the suffix j is $m+1$ or more, the processing is ignored.

8. An apparatus for manifesting characteristic existing in a symbolic sequence I_j ($j = 1 \sim m$), comprising:

a means for memorizing the symbolic sequence I_j ;

a means for effecting conversion to a parallel sequence $A(k)$ of partial symbolic sequences in which the suffix j is aligned in the following positional relation:

$j = 1, 2, \dots, k-1, k$

$j = k+1, k+2, \dots, k+k-1, k+k$

:

:

:

$J = (n-1)k+1, (n-1)k+2, \dots, (n-1)k+k-1, (n-1)k+k$

$j = nk+1, nk+2, \dots, nk+k-1, nk+k$

or

$j = 1, 2, \dots, k-1, k$

$j = k+k, k+k-1, \dots, k+2, k+1$

:

:

:

$j = (n-1)k+k, (n-1)k+k-1, \dots, (n-1)k+2, (n-1)k+1$

$j = nk+1, nk+2, \dots, nk+k-1, nk+k;$

and

a means for outputting the converted parallel sequence $A(k)$ using one or more expression means selected from hue, lightness and saturation of color and from interval, tone and volume of sound;

wherein, k represents an integer of 2 or more, and n represents an integer such that $nk < m \leq nk+k$, and when the suffix j is $m+1$ or more, the processing is ignored.

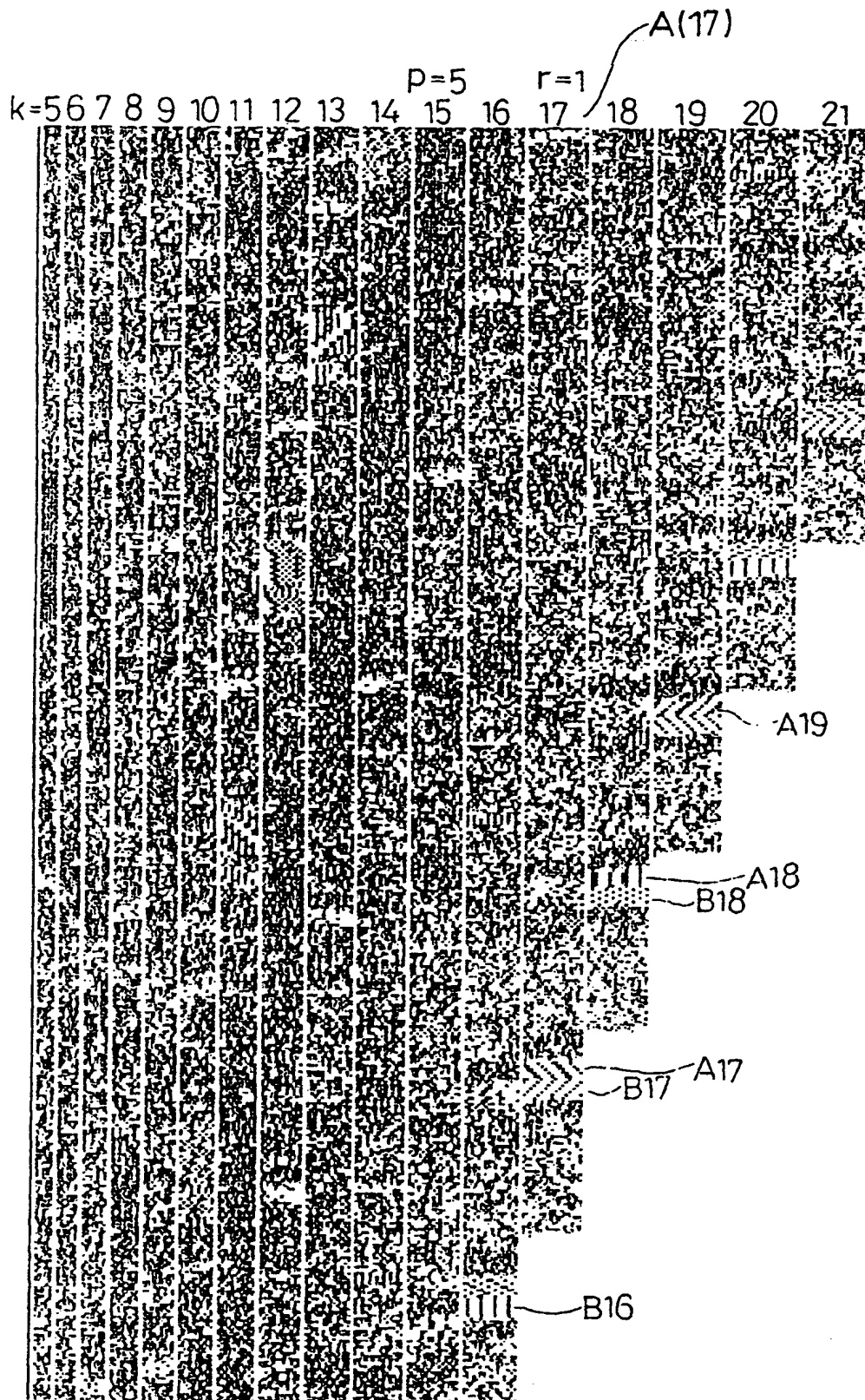


FIG.1

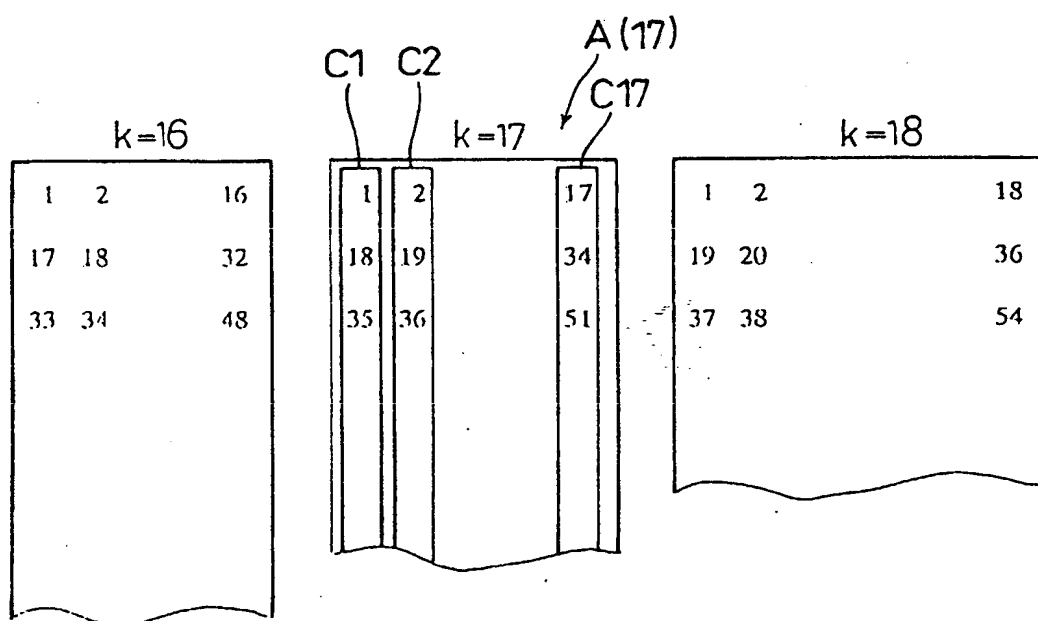


FIG.2

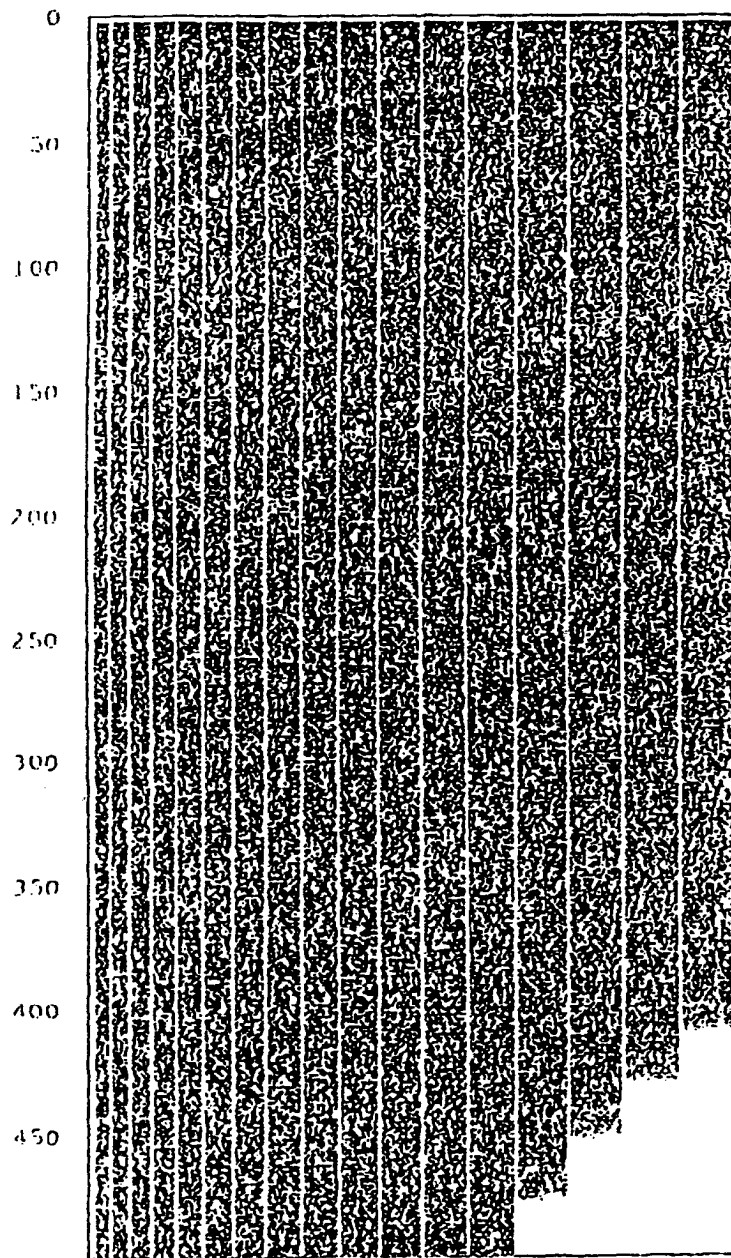


FIG.3

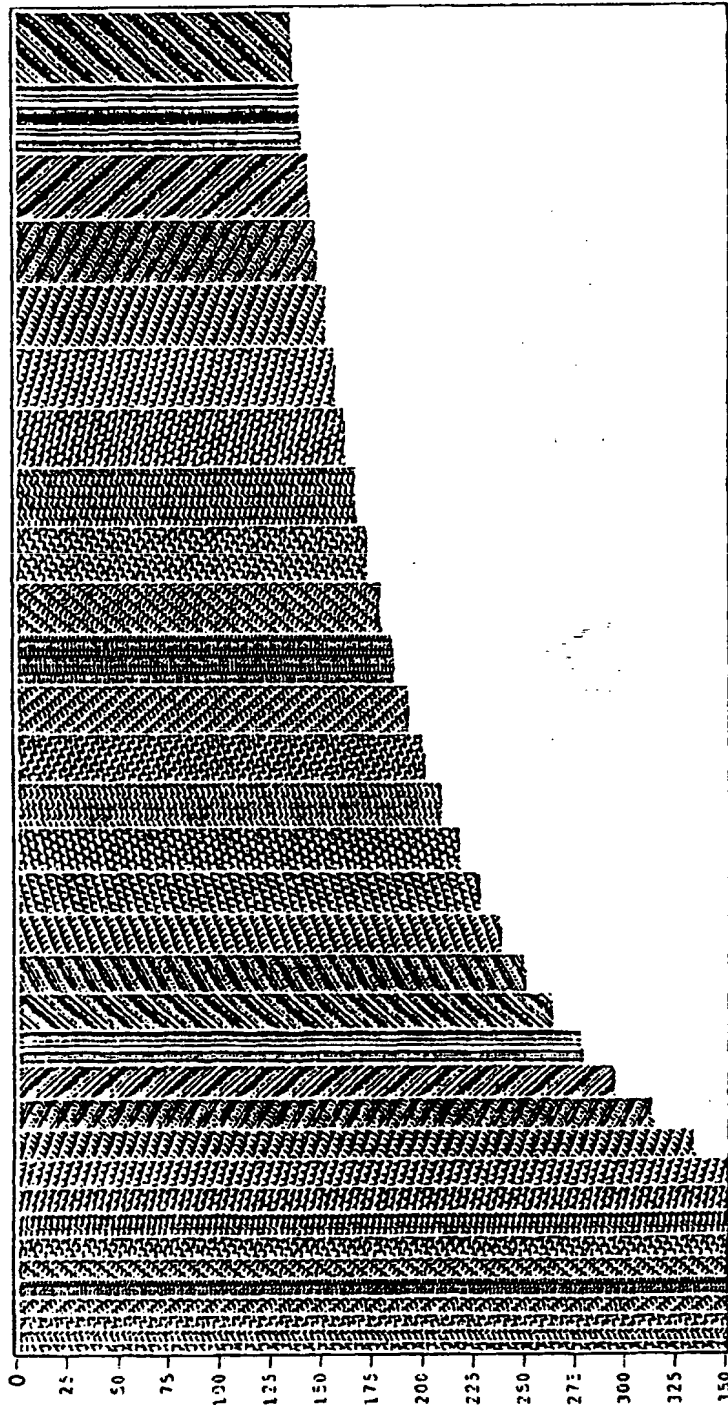


FIG.4

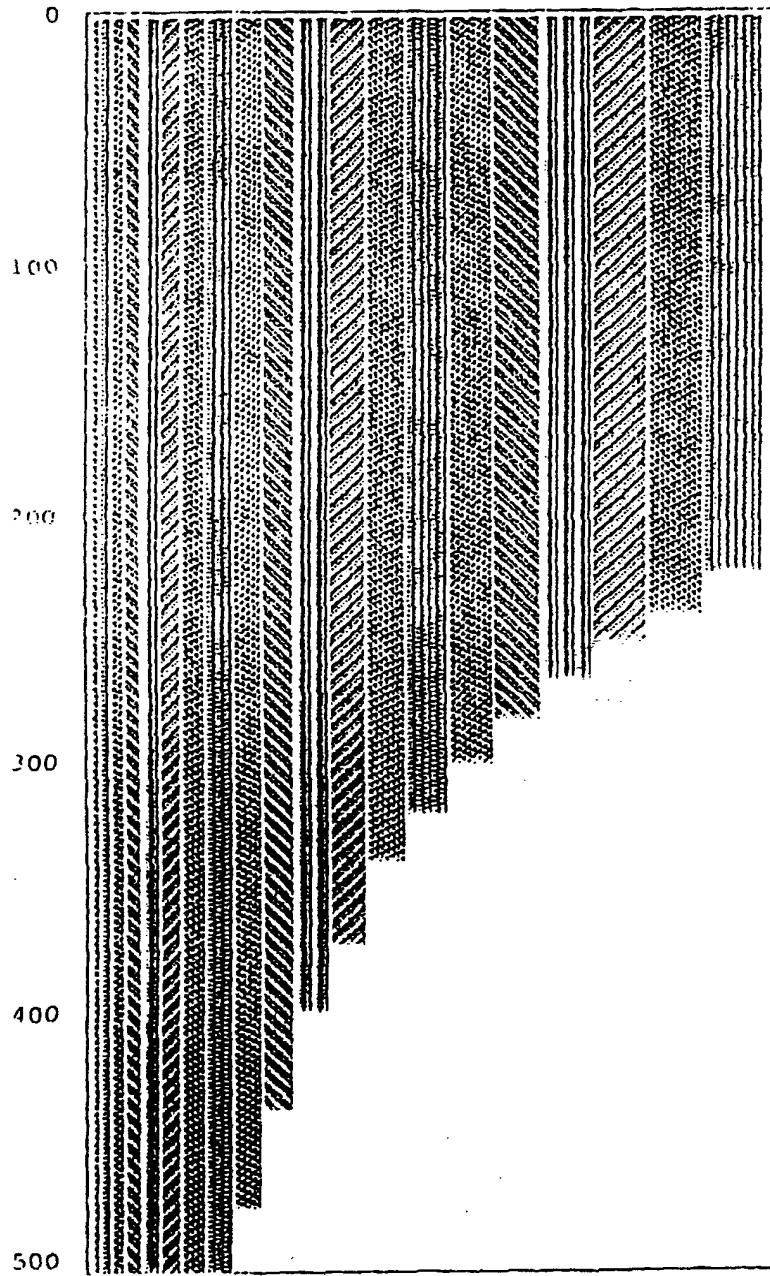
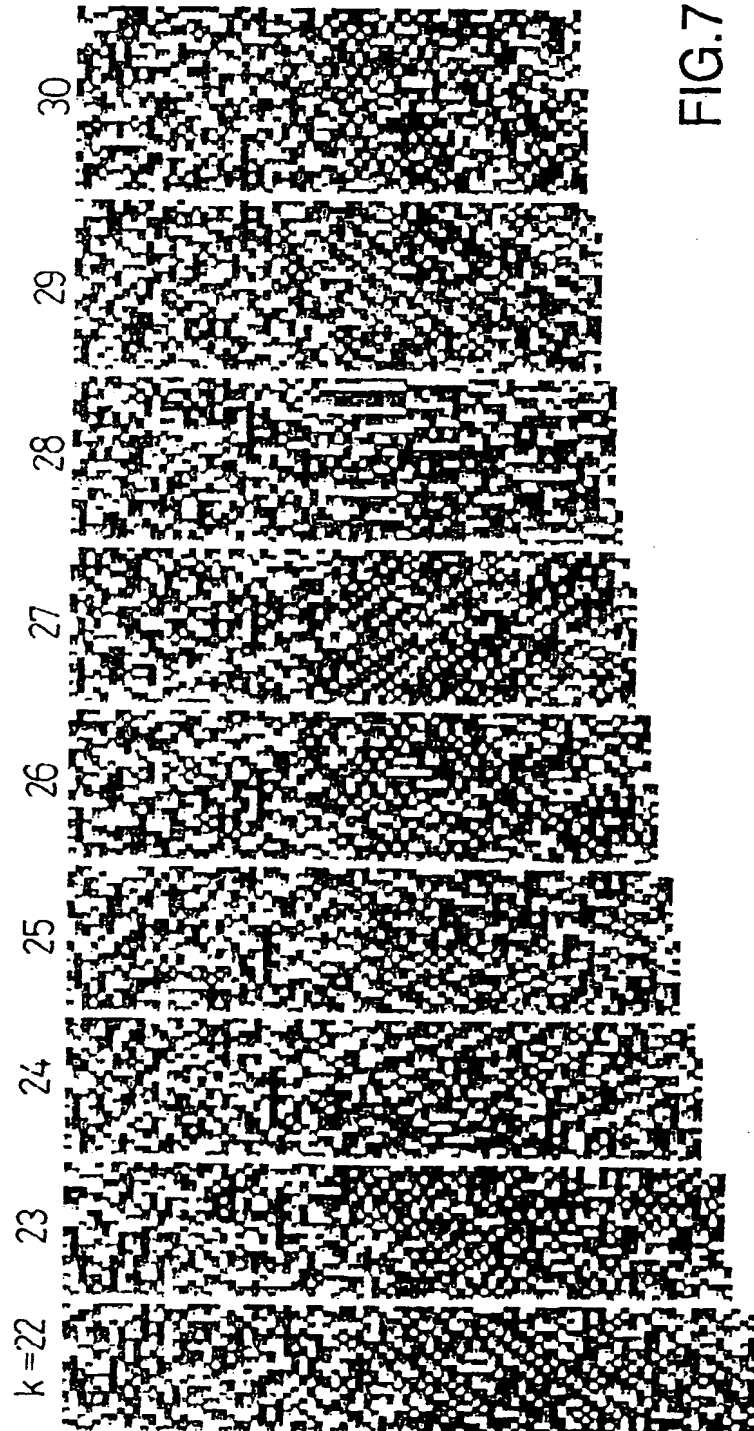


FIG.5



FIG.6



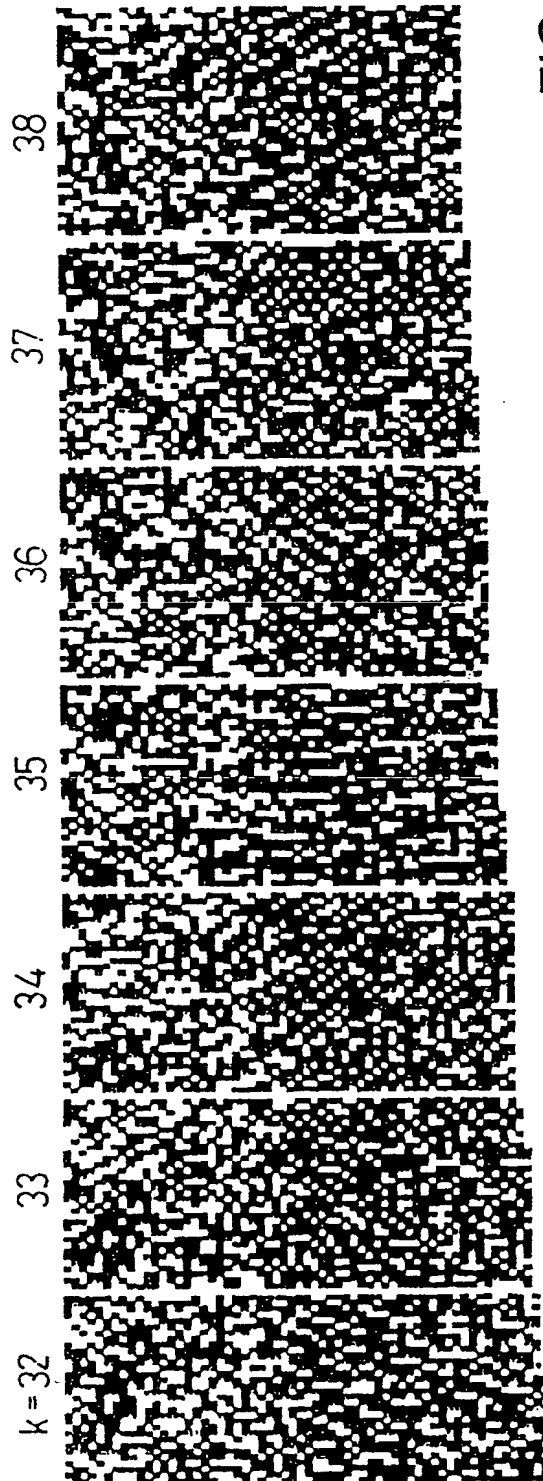


FIG.8

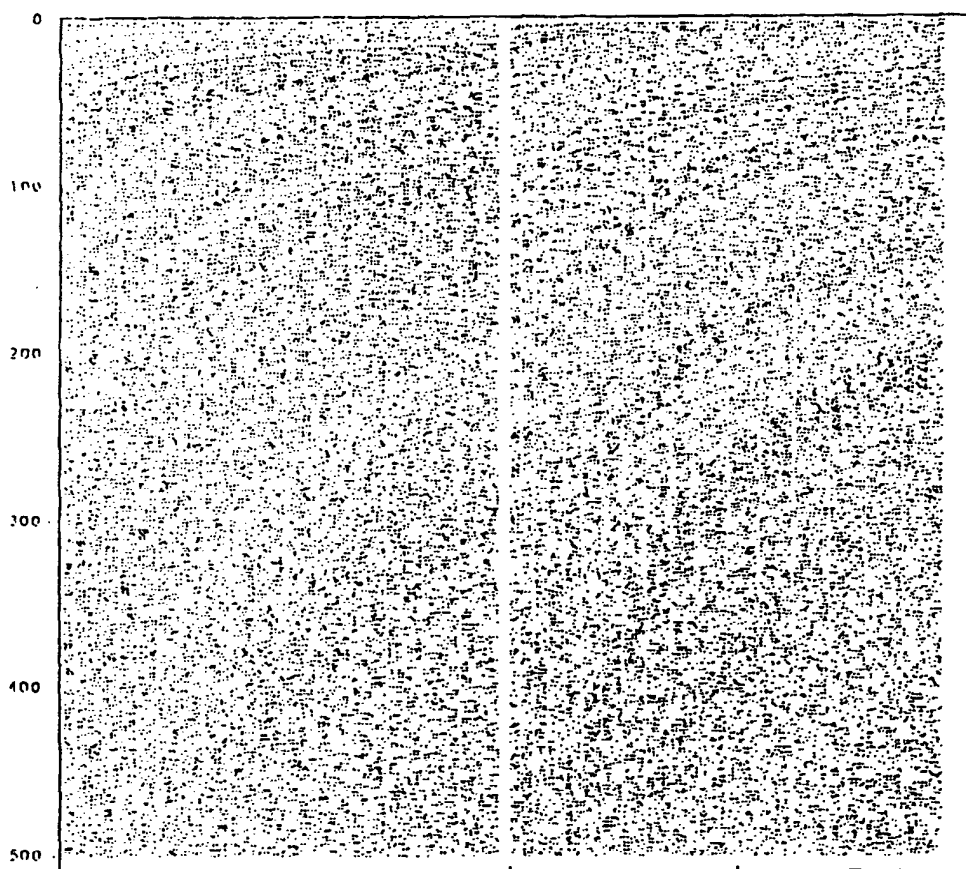


FIG.9

	k=2	k=3	k=4	k=20
1	1 2	1 2 3	1 2 3 4	1 2 3 . . . 20
2	3 4	4 5 6	5 6 7 8	21
3	5 6	7 8 9	9 10 11 12	. . .
4	7 8	10 11 12	13 14 15 16	
5	9 10	13 14 15	17 18 19 20 100
6	11 12	16 17 18	
7	13 14	19 20	
8	15 16	
9	17 18	. . .		
10	19 20	. . .		
11	.			
12	.			
13	.			
14				
15				
16				
17				
18				
19				
20				
.				
.				
.				

FIG.10

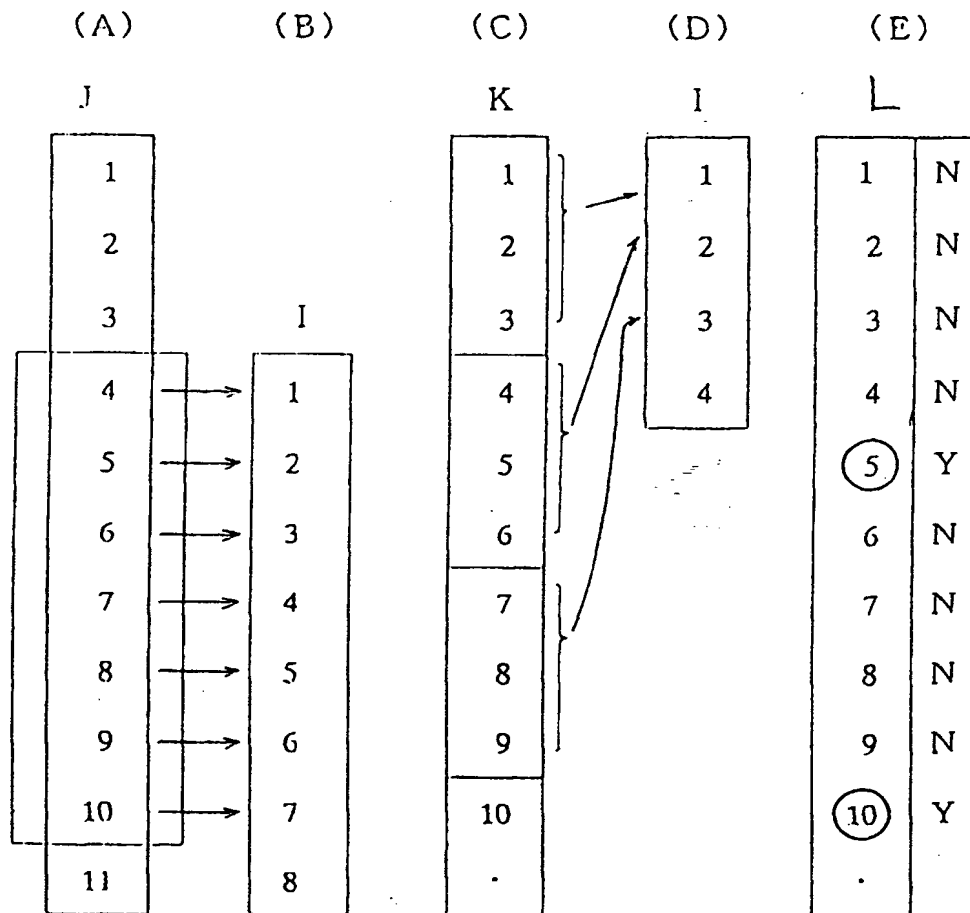


FIG.11

	$q = 5$ $k = 2$	$k = 3$	$k = 4$	$k = 6$
1	5 10	5 10 15	5 10 15 20	5 10 15 20 25 30
2	15 20	20 25 30	25 30 35 40	35 40 45 50 55 60
3	25 30	35 40 45	45 50 55 60	65 70 75 80 85 90
4	35 40	50 55 60	65 70 75 80	95
5	45 50	65 70 75	85 90 95	
6	55 60	80 85 90		
7	65 70	95		
8	75 80			
9	85 90			
10	95			
11				
12				
13				
14				
15				
16				
17				
18				
19				
20				
.				
.				
.				

FIG.12

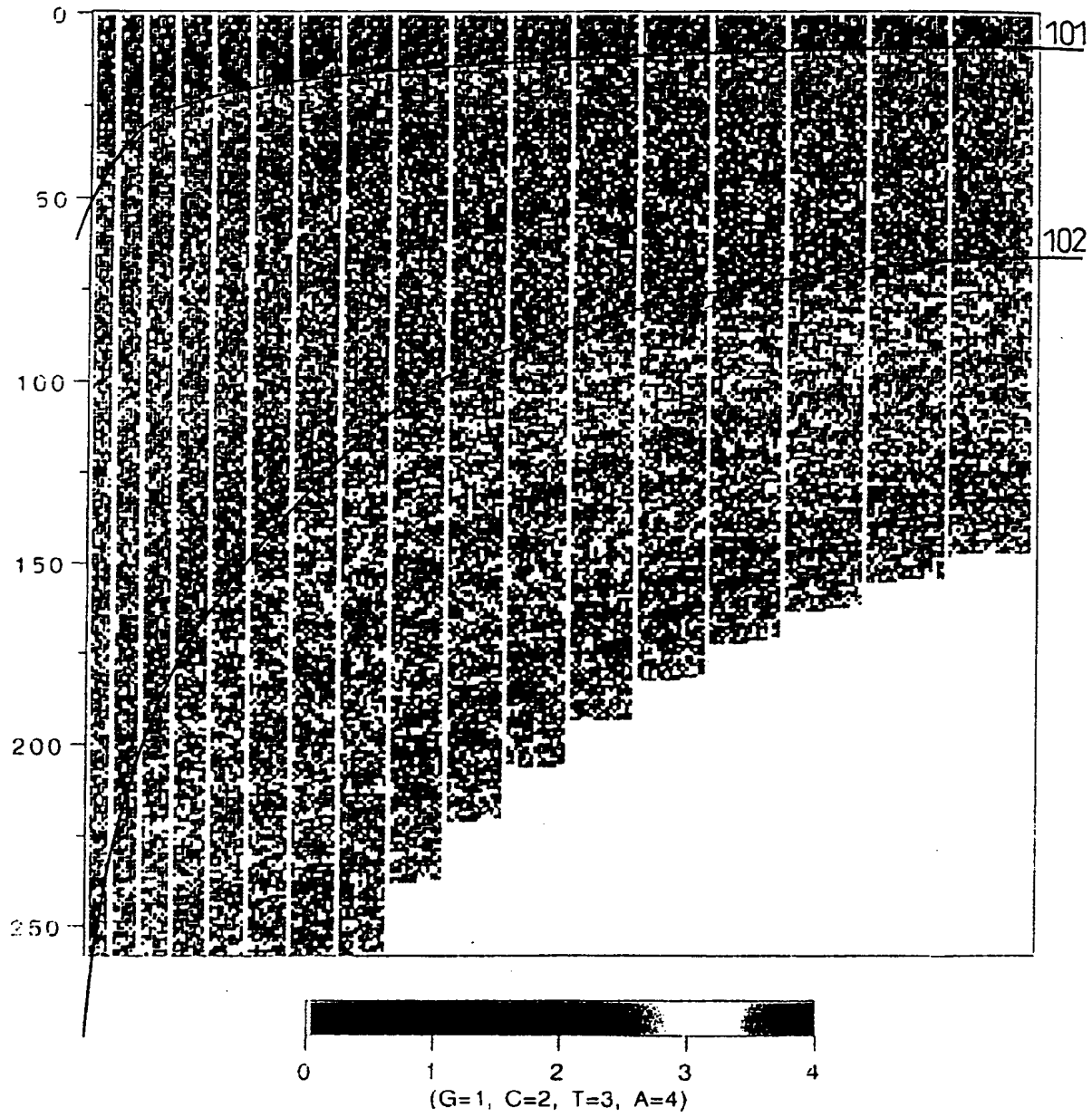


FIG.13

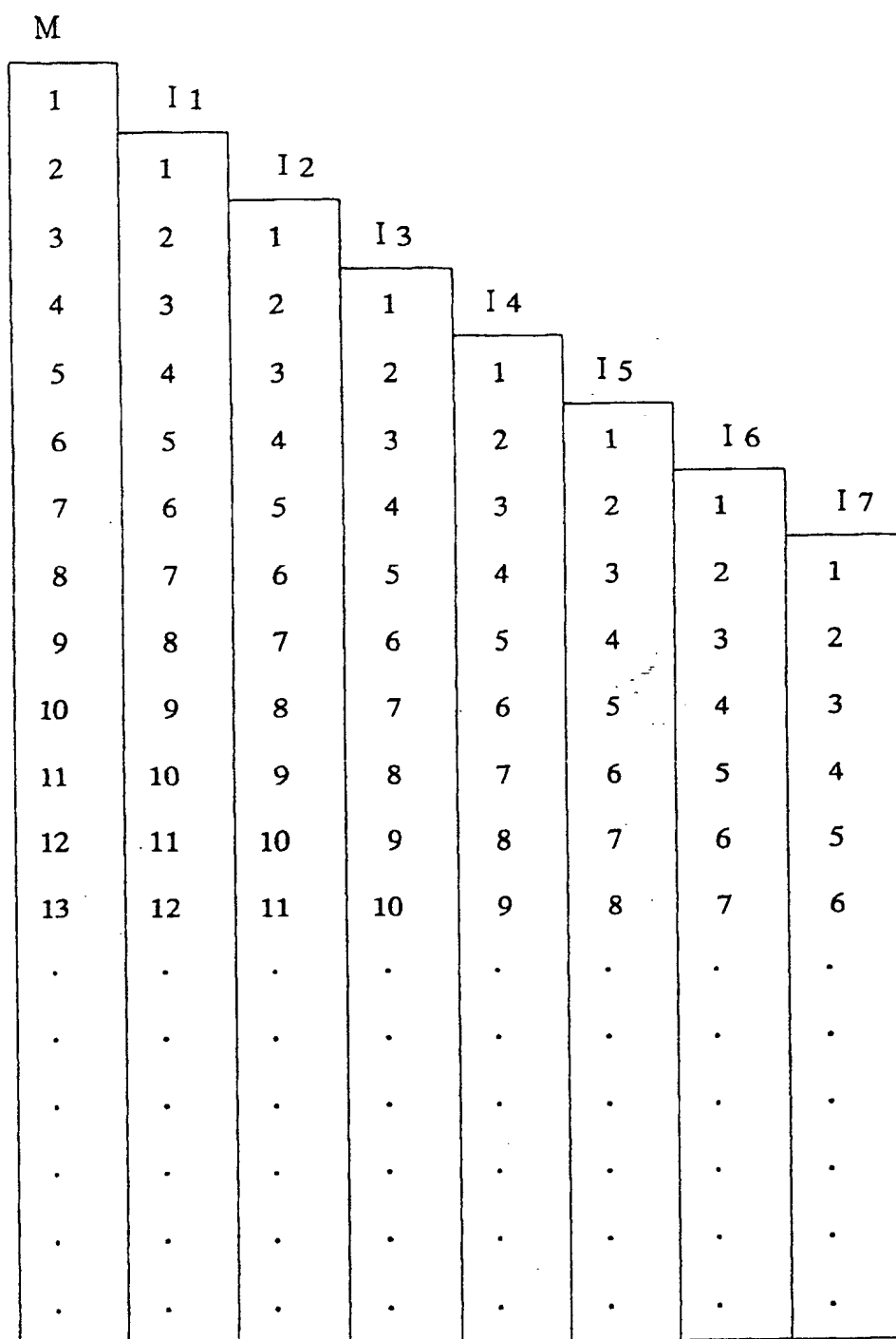
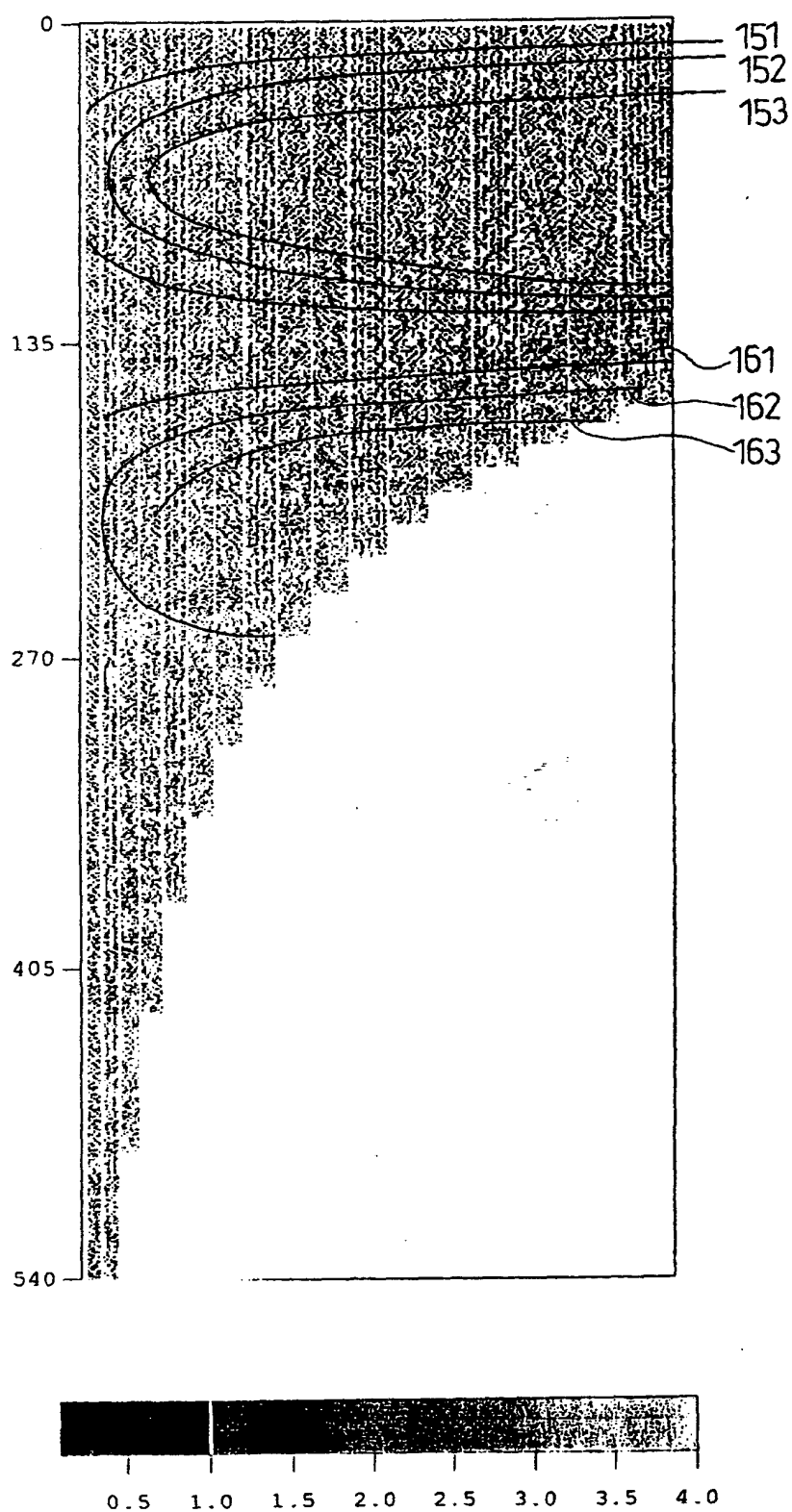


FIG.14



G=1, C=2, T=3, A=4

FIG.15

	k=2	k=3	k=4	k=20
1	1 2	1 2 3	1 2 3 4	1 2 3 . . . 20
2	4 3	6 5 4	8 7 6 5	. . . 22 21
3	5 6	7 8 9	9 10 11 12	
4	8 7	12 11 10	16 15 14 13	
5	9 10	13 14 15	17 18 19 20 100
6	12 11	18 17 16	
7	13 14	19 20	
8	16 15	
9	17 18	. . .		
10	20 19	. . .		
11	.			
12	.			
13	.			
14				
15				
16				
17				
18				
19				
20				
.				
.				
.				

FIG.16

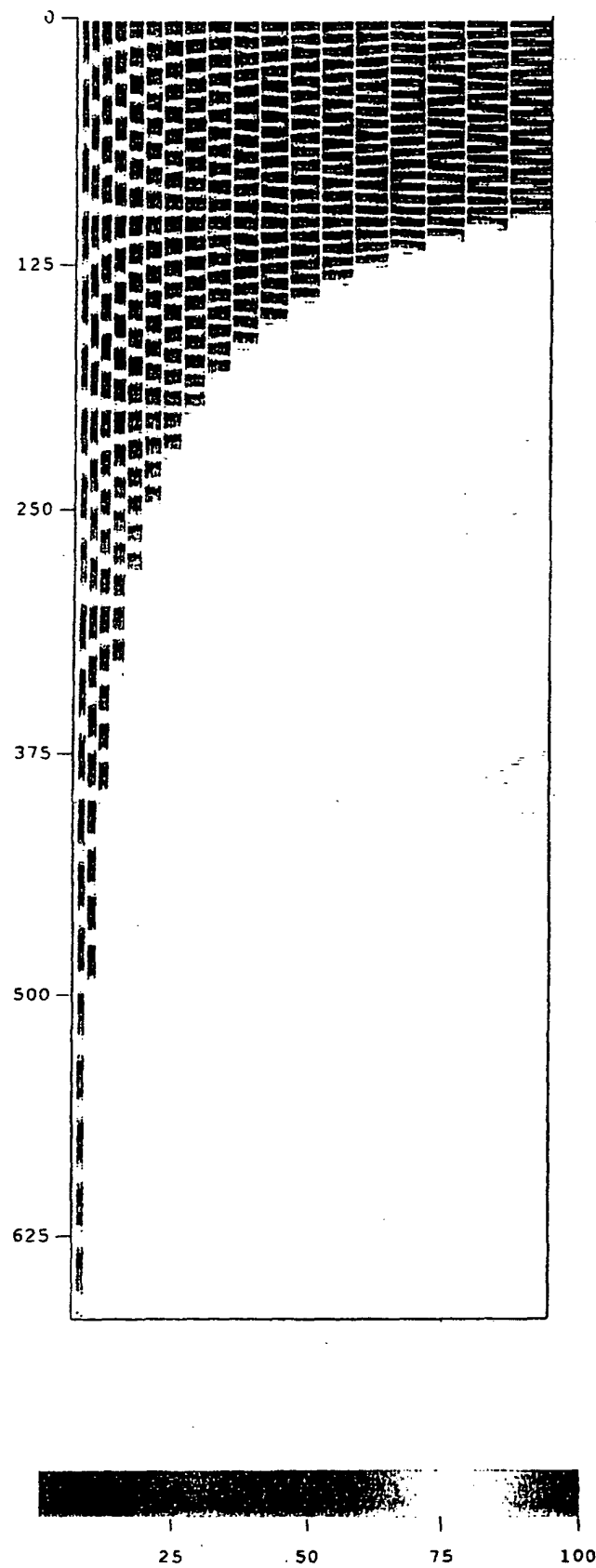
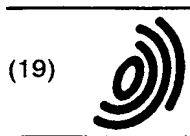


FIG.17



FIG.18



Europäisches Patentamt
European Patent Office
Offic européen d s br vets



(11) **EP 0 898 236 A3**

(12) **EUROPEAN PATENT APPLICATION**

(88) Date of publication A3:
10.01.2001 Bulletin 2001/02

(51) Int. Cl.⁷: **G06F 17/30, G06F 19/00**

(43) Date of publication A2:
24.02.1999 Bulletin 1999/08

(21) Application number: **98115643.3**

(22) Date of filing: **19.08.1998**

(84) Designated Contracting States:
**AT BE CH CY DE DK ES FI FR GB GR IE IT LI LU
MC NL PT SE**
Designated Extension States:
AL LT LV MK RO SI

(30) Priority: **20.08.1997 JP 22390897**

(71) Applicant:
**Toa Gosei Kabushiki Kaisha
Tokyo 105-0003 (JP)**

(72) Inventors:
• **Yoshida, Tetsuhiko**
Nagoya-shi, Aichi 455-0022 (JP)
• **Oosawa, Kenji**
Nara-shi, Nara 631-0011 (JP)
• **Obata, Nobuaki**
Nagoya-shi, Aichi 464-0096 (JP)

(74) Representative: **Kuhnen & Wacker**
Patentanwalts-gesellschaft mbH,
Alois-Steinecker-Strasse 22
85354 Freising (DE)

(54) **Method and apparatus for manifesting characteristic existing in symbolic sequence**

(57) A method which manifests characteristic which is latent and can not be recognized, although it exists in a complicated symbolic sequence, for example, a nucleotide sequence of DNA, and thereby enables recognition of the characteristic unrecognized yet, is provided.

When a symbolic sequence I_j ($j = 1 \sim m$) is given, there is an effected conversion to a parallel sequence $A(k)$ of partial symbolic sequences in which the suffix j is aligned in the following positional relation:

$j =$	$1,$	$2, \dots$	$k-1,$	k
$j =$	$k+1,$	$k+2, \dots$	$k+k-1,$	$k+k$
:				
:				
$J =$	$(n-1)k+1,$	$(n-1)k+2, \dots$	$(n-1)k+k-1,$	$(n-1)k+k$
$j =$	$nk+1,$	$nk+2, \dots$	$nk+k-1,$	$nk+k$

and $A(k)$ is formed with changing k to $p, p+r, p+2r, p+3r, \dots$, and the whole parallel sequences $\Sigma A(k)$ is obtained.

When regularity of period length k exists in the symbolic sequence I_j , the regularity remarkably appears in the partial symbolic sequences obtained by extracting one symbol at every $k-1$ symbols from the symbolic sequence.

EP 0 898 236 A3

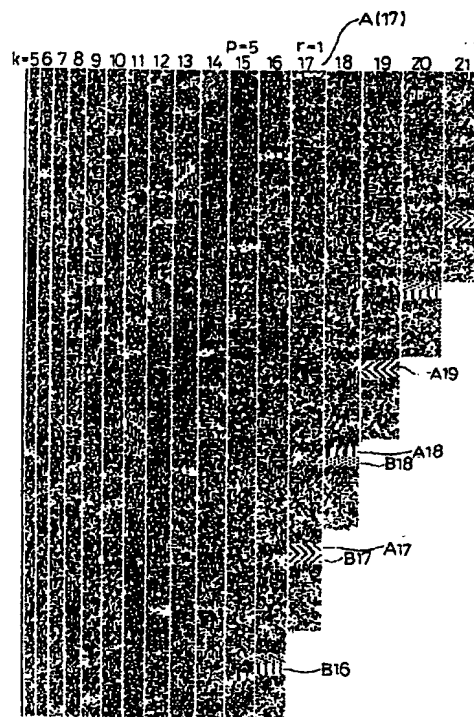


FIG.1



European Patent
Office

EUROPEAN SEARCH REPORT

Application Number
EP 98 11 5643

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (Int.CI.8)
P,X	OBATA N ET AL: "A method for information analysis of sequences by two-dimensional pattern formation with coloration" PROCEEDINGS OF THE FIRST INTERNATIONAL CONFERENCE. QUANTUM INFORMATION, PROCEEDINGS OF FIRST INTERNATIONAL CONFERENCE ON QUANTUM INFORMATION, NAGOYA, JAPAN, 4-8 NOV. 1997, pages 27-58, XP000961400 1999, Singapore, World Scientific, Singapore ISBN: 981-02-3934-3 * page 31, paragraph 2.1 - page 33, paragraph 2.2 *	1-8	G06F17/30 G06F19/00
A	MAIZEL J V ET AL: "ENHANCED GRAPHIC MATRIX ANALYSIS OF NUCLEIC ACID AND PROTEIN SEQUENCES" PROCEEDINGS OF THE NATIONAL ACADEMY OF SCIENCES OF THE USA,US,NEW YORK, NY, vol. 78, no. 12, 1 December 1981 (1981-12-01), pages 7665-7669, XP002036657 * page 7666, right-hand column, line 57-62 * * page 7666, right-hand column, line 19-21 *	1-8	TECHNICAL FIELDS SEARCHED (Int.CI.6) G06F
A	EP 0 561 563 A (AMERICAN TELEPHONE & TELEGRAPH) 22 September 1993 (1993-09-22) * page 2, line 45 - page 3, line 9 * -/--	1-8	
The present search report has been drawn up for all claims			
Place of search THE HAGUE		Date of completion of the search 6 November 2000	Examiner Correia Martins, F
<p>CATEGORY OF CITED DOCUMENTS</p> <p>X : particularly relevant if taken alone Y : particularly relevant if combined with another document of the same category A : technological background O : non-written disclosure P : intermediate document</p> <p>T : theory or principle underlying the invention E : earlier patent document, but published on, or after the filing date D : document cited in the application L : document cited for other reasons & : member of the same patent family, corresponding document</p>			

EPO FORM 1503 03 83 B2 (P04C01)



European Patent
Office

EUROPEAN SEARCH REPORT

Application Number
EP 98 11 5643

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (Int.Cl.6)
A	<p>X. GUAN ET. AL.: "A FAST LOOK-UP ALGORITHM FOR DETECTING REPETITIVE DNA SEQUENCES"</p> <p>6 January 1996 (1996-01-06) , WORLD SCIENTIFIC , SAN FRANCISCO, USA</p> <p>XP000953382</p> <p>* the whole document *</p> <p>-----</p>	1-8	
			TECHNICAL FIELDS SEARCHED (Int.Cl.6)
The present search report has been drawn up for all claims			
Place of search		Date of completion of the search	Examiner
THE HAGUE		6 November 2000	Correia Martins, F
<p>CATEGORY OF CITED DOCUMENTS</p> <p>X : particularly relevant if taken alone</p> <p>Y : particularly relevant if combined with another document of the same category</p> <p>A : technological background</p> <p>O : non-written disclosure</p> <p>P : intermediate document</p> <p>T : theory or principle underlying the invention</p> <p>E : earlier patent document, but published on, or after the filing date</p> <p>D : document cited in the application</p> <p>L : document cited for other reasons</p> <p>& : member of the same patent family, corresponding document</p>			

EPO FORM 1503 03/82 (P4C01)

**ANNEX TO THE EUROPEAN SEARCH REPORT
ON EUROPEAN PATENT APPLICATION NO.**

EP 98 11 5643

This annex lists the patent family members relating to the patent documents cited in the above-mentioned European search report.
The members are as contained in the European Patent Office EDP file on
The European Patent Office is in no way liable for these particulars which are merely given for the purpose of information.

06-11-2000

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
EP 0561563 A	22-09-1993	US 5953006 A	14-09-1999
		CA 2091503 A,C	19-09-1993
		JP 6043838 A	18-02-1994
		US 5627748 A	06-05-1997
		US 5511159 A	23-04-1996
<hr/>			

